

Principal component analysis

Ethan Y. Jaffe

Consider n vectors $x_1, \dots, x_n \in \mathbf{R}^N$. One often wishes to find the best lower-dimensional representation of the vectors x_1, \dots, x_n , i.e. to find the closest k -dimensional affine hyperplane ($k < N$) to the vectors. In other words, one wishes to minimize, over k -dimensional hyperplanes (through the origin) and offsets $b \in \mathbf{R}^N$ the quantity

$$\sum_{i=1}^n \text{dist}(x_i, V + b)^2 = \sum_{i=1}^n \|(1 - \Pi_V)(x_i - b)\|^2, \quad (1)$$

where Π_V denotes the orthogonal projection onto the subspace V .

Principal component analysis (PCA) provides an exact answer to this problem in terms of the eigenvectors of the sample covariance matrix¹ $\Sigma = XX^T$, where X is the matrix whose columns are the vectors $x_i - \hat{x}$, \hat{x} denoting the mean $\hat{x} = \frac{x_1 + \dots + x_n}{n}$. Since Σ is positive semi-definite, it has eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N \geq 0$, and corresponding orthonormal eigenvectors v_1, \dots, v_N .

Theorem 1 (PCA). *Quantity (1) is minimized if and only if V is the span of k orthonormal eigenvectors of Σ corresponding to the k largest eigenvalues of Σ , and $b - \hat{x} \in V$. In this case, the minimum distance is $\sum_{i=1}^n \|x_i\|^2 - \left(\sum_{j=1}^k \lambda_j\right)$.*

Remark 2. Notice that this is an if and only if statement. It asserts how to find a pair (V, b) minimizing (1) and also characterizes the form of any minimizing pair.

Remark 3. If $k \geq 1$, then the offset b is only unique up to an element of V . This is to be expected, since there are many pairs (V, b) defining the same affine hyperplane. The hyperplane V also need not be unique, because the multiplicity of the eigenvalues λ_i is not necessarily 1.

Remark 4. Another way of interpreting the theorem is that PCA finds the directions of maximum variance. The PCA procedure may thus be explained as follows: demean the data, and then extract a set of k orthogonal directions in which the variance of the data (i.e. the vectors x_1, \dots, x_n) is maximal.

We will prove the theorem in three steps. In the first, we assert that the minimum is actually realized at a pair (V, b) , and that for this pair $b - \hat{x} \in V$. This allows us to reduce

¹Really only the sample covariance matrix up to a constant factor depending on n ; this does not affect any of the statements made below.

the case that the mean $\hat{x} = 0$ and minimizing pairs (V, b) where $b = 0$. In the second step we relate this minimization problem to a maximization problem involving Σ in a purely linear-algebraic fashion. In the third step, we solve this maximization problem.

Lemma 5. *The minimum to (1) is attained a pair (V, b) . Furthermore, if (V, b) is any minimal pair, then $b - \hat{x} \in V$.*

Proof. Let us start with the proof of the second statement. Notice that $(1 - \Pi_V)\hat{x}$ is the mean of $(1 - \Pi_V)x_1, \dots, (1 - \Pi_V)x_n$, and so

$$\sum_{i=1}^n \|(1 - \Pi_V)x_i - (1 - \Pi_V)\hat{x}\|^2 \leq \sum_{i=1}^n \|(1 - \Pi_V)x_i - c\|^2, \quad (2)$$

for any $c \in \mathbf{R}^N$, with equality holding if and only if $c = (1 - \Pi_V)\hat{x}$. In particular this is true for $c = (1 - \Pi_V)b$, which means that if (V, b) is a minimiaing pair, then $(1 - \Pi_V)(b - \hat{x}) = 0$, or $b - \hat{x} \in V$.

Let us now prove that the minimum is attained. Let (V_j, b_j) be a sequence of ordered pairs for which

$$\sum_{i=1}^n \|(1 - \Pi_{V_j})(x_i - b_j)\|^2 \rightarrow \inf_{(V,b)} \sum_{i=1}^n \|(1 - \Pi_V)(x_i - b)\|^2.$$

Using (2), we may assume that $(1 - \Pi_{V_j})b_j = (1 - \Pi_{V_j})\hat{x}$, or just that $b_j = \hat{x}$. Passing to a subsequence V_{j_ℓ} , we may assume that there exists a hyperplane V for which $\Pi_{V_{j_\ell}} \rightarrow \Pi_V$. Indeed, letting $(v_{j,1}, \dots, v_{j,k})$ be an orthonormal basis of V_j , we may extract a convergent subsequence $(v_{j_\ell,1}, \dots, v_{j_\ell,k}) \rightarrow (v_1, \dots, v_k)$, an orthonormal basis of a hyperplane V . It is easy to check the convergence $\Pi_{V_{j_\ell}} \rightarrow \Pi_V$. Thus, $(1 - \Pi_{V_{j_\ell}})(x_i - b_j) = (1 - \Pi_{V_{j_\ell}})(x_i - \hat{x}) \rightarrow (1 - \Pi_V)(x_i - \hat{x})$, which means that

$$\sum_{i=1}^n \|(1 - \Pi_{V_{j_\ell}})(x_i - b_j)\|^2 \rightarrow \sum_{i=1}^n \|(1 - \Pi_V)(x_i - \hat{x})\|^2 = \inf_{(V,b)} \sum_{i=1}^n \|(1 - \Pi_V)(x_i - b)\|^2,$$

so that the mimium is achieved at (V, \hat{x}) . □

Given this lemma, we see that the minimization problem is equivalent to minimizing with $b = \hat{x}$ fixed, and that if (V, b') is a minimizing pair then (V, \hat{x}) is a minimizing pair, too, since $b' - \hat{x} \in V$, and so (1) is unaffected by replacing b' with \hat{x} . Thus, without loss of generality, we may assume that $\hat{x} = 0$, and focus on the minimization problem without offset, i.e. with $b = 0$.

Lemma 6. *The following identity holds for any k -dimensional hyperplane V :*

$$\sum_{i=1}^n \|(1 - \Pi_V)x_i\|^2 = \sum_{i=1}^n \|x_i\|^2 - \text{Tr}(\Pi_V \Sigma).$$

In particular, the left-hand side is minimized when $\text{Tr}(\Pi_V \Sigma)$ is maximized.

Proof. First notice that

$$\sum_{i=1}^n \|(1 - \Pi_V)x_i\|^2 = \sum_{i=1}^n \|x_i\|^2 - \sum_{i=1}^n \|\Pi_V x_i\|^2,$$

so we only need to show that

$$\sum_{i=1}^n \|\Pi_V x_i\|^2 = \text{Tr}(\Pi_V \Sigma).$$

Let v_1, \dots, v_k be an orthonormal basis of V , and extend it via v_{k+1}, \dots, v_N to an orthonormal basis of \mathbf{R}^N . Then

$$\text{Tr}(\Pi_V \Sigma) = \sum_{j=1}^N \langle \Pi_V \Sigma v_j, v_j \rangle = \sum_{j=1}^k \langle \Pi_V \Sigma v_j, v_j \rangle. \quad (3)$$

Now notice that by definition

$$\Sigma = \sum_{i=1}^n x_i x_i^T,$$

or in more coordinate-free language

$$\Sigma = \sum_{i=1}^n \langle x_i, \cdot \rangle x_i.$$

Thus for $1 \leq j \leq k$,

$$\langle \Pi_V \Sigma v_j, v_j \rangle = \sum_{i=1}^n \langle (\Pi_V x_i) \langle x_i, v_j \rangle, v_j \rangle = \sum_{i=1}^n \langle x_i, v_j \rangle \langle \Pi_V x_i, v_j \rangle = \sum_{i=1}^n |\langle x_i, v_j \rangle|^2.$$

Plugging this into (3) yields

$$\text{Tr}(\Pi_V \Sigma) = \sum_{i=1}^n \sum_{j=1}^k |\langle x_i, v_j \rangle|^2.$$

We recognize the inner sum as $\|\Pi_V x_i\|^2$. Thus

$$\text{Tr}(\Pi_V \Sigma) = \sum_{i=1}^n \|\Pi_V x_i\|^2,$$

which completes the proof. □

The crux of the proof is now to show the following general statement in linear algebra: let T be any positive-semidefinite $N \times N$ matrix, with eigenvalues $\mu_1 \geq \dots \geq \mu_N \geq 0$.

Proposition 7. *The quantity $\text{Tr}(\Pi_V T)$ is maximized over k -dimensional hyperplanes V if and only if V is the span of k orthonormal eigenvectors of T corresponding to its k largest eigenvalues. In this case, the maximum is $\sum_{j=1}^k \mu_j$.*

Proof. Let e_1, \dots, e_N be an orthonormal basis of eigenvectors of T , corresponding (in order) to the eigenvalues $\mu_1 \geq \dots \geq \mu_N$ of T . Let V be any k -dimensional hyperplane. Then

$$\text{Tr}(\Pi_V T) = \sum_{j=1}^N \langle \Pi_V T e_j, e_j \rangle = \sum_{j=1}^N \mu_j \langle \Pi_V e_j, e_j \rangle = \sum_{j=1}^N \mu_j \|\Pi_V e_j\|^2. \quad (4)$$

Notice that if $V = \text{span}\{e_1, \dots, e_k\}$, then this quantity is precisely $\sum_{j=1}^k \mu_j$. To prove that this is the maximum quantity, we need only bound (4). Set $a_j = \|\Pi_V e_j\|^2$. Then for all $1 \leq j \leq N$, $0 \leq a_j \leq 1$. Furthermore, $\sum_{j=1}^N a_j = \sum_{j=1}^N \langle \Pi_V e_j, e_j \rangle = \text{Tr}(\Pi_V) = k$. Thus, from (4),

$$\begin{aligned} \text{Tr}(\Pi_V T) &= \sum_{j=1}^N a_j \mu_j \\ &= \sum_{j=1}^k \mu_j + \sum_{j=1}^k \mu_j (a_j - 1) + \sum_{j=k+1}^N a_j \mu_j \\ &\leq \sum_{j=1}^k \mu_j + \sum_{j=1}^k \mu_k (a_j - 1) + \sum_{j=k+1}^N a_j \mu_k \\ &= \sum_{j=1}^k \mu_j + \mu_k \sum_{j=1}^N a_j - \mu_k \sum_{j=1}^k 1 \\ &= \sum_{j=1}^k \mu_j + k\mu_k - k\mu_k = \sum_{j=1}^k \mu_j. \end{aligned} \quad (5)$$

This shows that the maximum is $\sum_{j=1}^k \mu_j$ and is achieved if V is the span of k orthonormal eigenvectors of T corresponding to its k largest eigenvalues.

We are not quite done, since we also need to show that the maximum is attained only if V is the span of k orthonormal eigenvectors of T corresponding to its k largest eigenvalues. This would be true, for instance if $\Pi_V e_j = e_j$ if $j \leq k$ and $\Pi_V e_j = 0$ if $j > k$, as in this case $V = \text{span}\{e_1, \dots, e_k\}$. Unfortunately, this need not be the case, because the multiplicity of μ_k could be greater than 1, and hence the choice of eigenvectors may not be unique. Instead, we will show that there is another orthonormal basis e'_1, \dots, e'_N (with corresponding eigenvalues $\mu_1 \geq \dots \geq \mu_N \geq 0$) for which $\Pi_V e'_j = e'_j$ if $j \leq k$ and $\Pi_V e'_j = 0$ if $j > k$. To do so, we will need to examine the inequality (5) more carefully. If $\text{Tr}(\Pi_V T)$ attains its maximum value $\sum_{j=1}^k \mu_j$ at a hyperplane V , then the inequality (5) must be an equality. This means that for $j \leq k$ either $a_j = 1$ or else $\mu_j = \mu_k$ and likewise for $j > k$ either $a_j = 0$ or else

$\mu_j = \mu_k$. Let $k_1 \leq k$ be the minimum index for which $\mu_{k_1} = \mu_k$ and likewise $k_2 \geq k$ be the maximum index for which $\mu_{k_2} = \mu_k$. In other words²

$$\mu_1 \geq \cdots \mu_{k_1-1} > \mu_{k_1} = \cdots = \mu_k = \cdots = \mu_{k_2} > \mu_{k_2+1} \geq \cdots \mu_N \geq 0.$$

Let E_k denote the eigenspace of T corresponding to μ_k . In other words,

$$E_k = \text{span}\{e_{k_1}, \dots, e_k, \dots, e_{k_2}\}.$$

Rexpressing the previous conditions on the a_j , we have that $a_j = 1$ for $j < k_1$ and $a_j = 0$ for $j > k_2$. Thus $\Pi_V e_j = e_j$ for $j < k_1$ and $\Pi_V e_j = 0$ for $j > k_2$. It is not difficult to check that this means that $\Pi_V : E_k \rightarrow E_k$.³ Thus $\Pi_V|_{E_k}$ is a rank $k - (k_1 - 1)$ orthogonal projection. In particular, we may find orthonormal vectors $e'_{k_1}, \dots, e'_k \in E_k$ for which $\Pi_V e'_j = e'_j$ for $k_1 \leq j \leq k$ and $\Pi_V e'_j = 0$ for $k < j \leq k_2$. Setting $e'_j = e_j$ for $j' < k_1$ or $j' > k_2$, it follows that $\Pi_V e'_j = e_j$ for $j \leq k$ and $\Pi_V e'_j = 0$ for $j > k$. Thus we have found the desired orthonormal basis and hence completed the proof. \square

²Of course, k_1 may equal 1 and k_2 may equal N and so this expression is not completely rigorous; for instance if $k_1 = 1$, then it asserts that $\mu_1 > \mu_1$. However, it gets the point across sufficiently well.

³For instance by showing that if $k_1 \leq j \leq k_2$, and $j' < k_1$ or $j' > k_2$ that $\Pi_V e_j$ is orthogonal to $e_{j'}$. Indeed, $\langle \Pi_V e_j, e_{j'} \rangle = \langle e_j, \Pi_V e_{j'} \rangle = \delta \langle e_j, e_{j'} \rangle$, where δ is either 0 or 1. Either way the quantity is 0. Thus $\Pi_V E_k$ is orthogonal to the span of all eigenvectors for eigenvalues other than μ_k . By the spectral theorem, this means that $\Pi_V E_k \subseteq E_k$.